

Density Estimation via Adaptive Partition and Discrepancy Control

Kun Yang, Wing Hung Wong

Abstract—Given iid samples from some unknown continuous density on hyper-rectangle $[0, 1]^d$, we attempt to learn a piecewise constant function that approximates this underlying density nonparametrically. Our density estimate is defined on a binary split of $[0, 1]^d$ and built up sequentially according to discrepancy criteria; the key ingredient is to control the discrepancy adaptively in each sub-rectangle to achieve overall bound. We prove that the estimate, even though simple as it appears, preserves most of the estimation power. By exploiting its structure, it can be directly applied to some important pattern recognition tasks such as mode seeking and density landscape exploration, we demonstrate its applicability through simulations and examples.

Index Terms—Binary Partition, Star Discrepancy, Density Estimation, Mode Seeking, Level Set Tree



1 INTRODUCTION

Classic empirical distribution (ED) and kernel density estimation (KDE) play an important role in nonparametric density estimation. Besides their long noticed drawbacks (e.g., ED is non-continuous; KDE is sensitive to the choice of bandwidth and scales poorly in high dimensions), they are not good summarization tools in dealing with data with high dimension and large size, e.g., evaluating them involves each data point and their functional forms provide little direct information of the “landscape” of the distribution. Density estimation based on domain partition dates back to the use of histogram, which is still an ubiquitous tool in data analysis today; however, its non-scalability in high dimensions limits its applications. Recently, Wong et al. [23] introduces optional pólya tree (OPT) as a class of conjugate nonparametric priors based on binary partition and optional stopping; later, Bayesian Sequential Partitioning (BSP) [15] is introduced as a computationally more attractive alternative to OPT and simulations show that the density constructed by BSP is very close to MAP estimate of OPT. The class of density functions obtained from them, which is piecewise constant on binary partitioned sample space, provide a concise summary of the data and can reveal the structure of the distribution in “multi-resolution”.

Motivated by previous work and the observation that the distribution conditioned on each sub-rectangle is uniform for piecewise constant densities, we construct the density estimator based on discrepancy criteria. We show that, in rather general setting, our estimated density, simple as it appears, preserves most of the estimation power.

Our algorithm, by exploiting the sequential build-up of binary partition as shown in Figure 1, can find the density efficiently.

It is important to distinguish our density estimator from OPT or BSP: 1) the density estimate of OPT or BSP is by sampling the posterior, while ours is constructed from a frequentist perspective and is deterministic; 2) the recursion in OPT has exponential complexity and BSP in principle searches among the exponential number of possible partitions, whereas our partitioning scheme is greedy and results in significant speedup; this point is illustrated through simulations, we also noticed that MAPs of OPT and BSP tend to overfit the data with noisy partitions, which raises difficulty in mode seeking; 3) as to binary partition, we no longer restrict the algorithm to split the hyper-rectangle evenly (in the middle); by introducing the “gap”, we do the partitioning more adaptive to the data; 4) OPT or BSP tries to control the estimation error directly; in contrast, our estimation is an indirect and weaker approximation to the true density that controls the integration error (under the same convergence rate as Monte Carlo methods) for the class of functions with finite total variation and finite variance.

2 DENSITY ESTIMATION VIA DISCREPANCY

Let Ω be a hyper-rectangle in \mathbb{R}^d . A binary partition \mathcal{B} on Ω is a collection of sub-rectangles whose union is Ω . Starting with $\mathcal{B}_1 = \{\Omega\}$ at level 1 and $\mathcal{B}_t = \{\Omega_1, \dots, \Omega_t\}$ at level t , \mathcal{B}_{t+1} is produced by dividing one of regions in \mathcal{B}_t along one of its coordinates, then combining both sub-rectangles with the rest of regions in \mathcal{B}_t ; continuing with this fashion, one can generate any binary partition at any level (Figure 1).

At each stage of sequential built-up of binary partition, to decide whether the sub-rectangle deserves further partitioning, we need to check whether the points in it

- Kun Yang is a PhD student in Institute of Computational and Mathematical Engineering, Stanford University, Stanford, CA 94305. E-mail: kunyang@stanford.edu.
- Wing Hung Wong is Stephen R. Pierce Family Goldman Sachs Professor in Science and Human Health, Professor of Statistics and Professor of Health Research and Policy at Stanford University, Stanford, CA 94305. E-mail: whwong@stanford.edu.

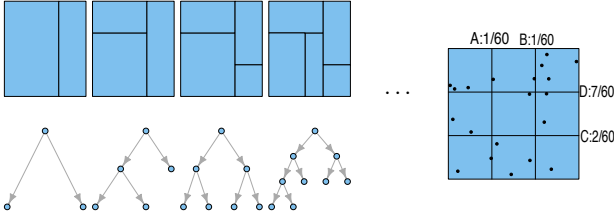


Fig. 1. Left: a sequence of binary partition and the corresponding tree representation; if we encode partitioning information (e.g., the location where the split occurs) in the nodes, the mapping is one-to-one. Right: the gaps with $m = 3$, we split the rectangle at location D, which corresponds to the largest gap, if it does not satisfy (5).

are “relative” uniformly scattered. Discrepancy, which is widely used in the analysis of Quasi-Monte Carlo methods, is a set of criteria to measure the uniformity of points in $[0, 1]^d$. The classic star discrepancy, which is used to bound the error of Quasi-Monte Carlo integration, is defined as,

Definition 2.1. The star discrepancy of $x_1, \dots, x_n \in [0, 1]^d$ is

$$D_n^*(x_1, \dots, x_n) = \sup_{a \in [0, 1]^d} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{x_i \in [0, a]} - \prod_{i=1}^d a_i \right| \quad (1)$$

The error bound is the famous Koksma-Hlawka inequality and the proof is included in [12],

Theorem 1. (Koksma-Hlawka inequality). Let $x_1, x_2, \dots, x_n \in [0, 1]^d$ and f be defined on $[0, 1]^d$, then

$$\left| \int_{[0, 1]^d} f(x) dx - \frac{1}{n} \sum_{i=1}^n f(x_i) \right| \leq D_n^*(x_1, \dots, x_n) V_{HK}^{[0, 1]^d}(f)$$

where $s = \{1, \dots, d\}$ and $V_{HK}^{[0, 1]^d}(f)$ is the total variation in the sense of Hardy and Krause, e.g., for any hyper-rectangle $[a, b]$, if all the involved partial derivatives of f are continuous on $[a, b]$, then

$$V_{HK}^{[a, b]}(f) = \sum_{u \subseteq \{1, \dots, d\}} \left\| \frac{\partial^{|u|} f}{\partial x_u} \Big|_{x_{s-u} = b_{s-u}} \right\|_1 \quad (2)$$

We split the sub-rectangle when the discrepancy of points in it is larger than some threshold value. In order to find a good location to split for $[a, b] = \prod_{j=1}^d [a_j, b_j]$, we divide j th dimension into m bins $[a_j, a_j + (b_j - a_j)/m, \dots, [a_j + (b_j - a_j)(m-2)/m, a_j + (b_j - a_j)(m-1)/m]$ and keep track of the gaps at $a_j + (b_j - a_j)/m, \dots, a_j + (b_j - a_j)(m-1)/m$, where the gap g_{jk} is defined as

$$g_{jk} = \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}(x_{ij} < a_j + (b_j - a_j) \frac{k}{m}) - \frac{k}{m} \right|$$

for $k = 1, \dots, (m-1)$, there are total $(m-1)d$ gaps recorded (Figure 1). $[a, b]$ is split into two sub-rectangles along the dimension and location corresponding to

maximum gap (Figure 1). The complete algorithm is given in Algorithm 1, and is explained in detail in the following sections.

The density, which is a piecewise constant function, is

$$\hat{p}(x) = \sum_{i=1}^n d(r_i) \mathbf{1}\{x \in r_i\} \quad (3)$$

where $\mathbf{1}$ is indicator function; $\{r_i, d(r_i)\}_{i=1}^n$ is a list of pairs of sub-rectangles and corresponding densities. Since the number of sub-regions is far less than data size, the partition is a concise representation of the data; in Experimental Results section, we demonstrate how $\hat{p}(x)$ can be leveraged in various machine learning applications.

Compared to histogram which has the same form (3) but suffers from curse of dimensionality, the rational behind our adaptive partition scheme is to avoid splitting the sub-rectangle where the data are relatively uniform. One classic results of histogram [9] states that if for each sphere S centered at the origin

$$\lim_{n \rightarrow \infty} \max_{r_i \cap S \neq \emptyset} \text{diam}(r_i) = 0$$

$$\lim_{n \rightarrow \infty} \frac{|\{r_i : r_i \cap S \neq \emptyset\}|}{n} = 0$$

then

$$\lim_{n \rightarrow \infty} \mathbf{E} \|p(x) - \hat{p}(x)\|_1 = 0$$

$$\lim_{n \rightarrow \infty} \|p(x) - \hat{p}(x)\|_1 = 0, a.s.$$

the key tool in proving its convergence is Lebesgue Density Theorem. However, our method can not guarantee the size of each sub-rectangle shrinks to 0, which causes the technical difficulty in proving its consistency. Instead, we establish a weaker convergence result in the following section and leave the pointwise convergence as an open problem.

3 THEORETICAL RESULTS

To establish our main theorem, we need the following three lemmas. Lemma 3.1 and 3.2 is trivial to show by (2) if f is smooth enough.

Lemma 3.1. Let f be defined on the hyper-rectangle $[a, b]$. Let $\{[a_i, b_i] : 1 \leq i \leq m < \infty\}$ be a split of $[a, b]$. Then

$$\sum_{i=1}^m V_{HK}^{[a_i, b_i]}(f) = V_{HK}^{[a, b]}(f)$$

Proof. The proof is in Lemma 1 of Section 5 in [17]. \square

Lemma 3.2. Let f be defined on the hyper-rectangle $[a, b]$. Let $\tilde{f}(\tilde{x})$ be defined on the hyper-rectangle $[\tilde{a}, \tilde{b}]$ by $\tilde{f}(\tilde{x}) = f(x)$ where $x_i = \phi_i(\tilde{x}_i)$ with ϕ_i is a strictly monotone (increasing or decreasing) invertible function from $[\tilde{a}_i, \tilde{b}_i]$ onto $[a_i, b_i]$, then

$$V_{HK}^{[\tilde{a}, \tilde{b}]}(\tilde{f}) = V_{HK}^{[a, b]}(f)$$

Algorithm 1 Density Estimation via Discrepancy (DED)

Let $P(\cdot)$ define the points and $\Pr(\cdot)$ define the probability mass in a hyper-rectangle respectively. W.L.G, we assume that $\Omega = [0, 1]^d$ and $P(\Omega) = \{x_i = (x_{i1}, \dots, x_{id}), x_i \in \Omega\}_{i=1}^{n_i}$ are iid samples drawn from an underlying distribution.

```

1: procedure DENSITY-ESTIMATOR( $\Omega, P, m, \theta$ )
2:    $\mathcal{B} = \{[0, 1]^d\}$ ,  $\Pr([0, 1]^d) = 1$ 
3:   while true do
4:      $\mathcal{B}' = \emptyset$ 
5:     for each  $r_i = [a_i, b_i]$  in  $\mathcal{B}$  do
6:       Calculate gaps  $\{g_{jk}\}_{j=1, \dots, d, k=1, \dots, m-1}$ 
7:       Scale  $P(r_i) = \{x_{ij}\}_{j=1}^{n_i}$  to  $\tilde{P} = \{\tilde{x}_{ij} = (\frac{x_{ij,1}-a_{i1}}{b_{i1}}, \dots, \frac{x_{ij,d}-a_{id}}{b_{id}})\}_{j=1}^{n_i}$ 
8:       if  $P(r_i) \neq \emptyset$  and  $D_{n_i}^*(\tilde{P}) > \alpha_i D_{n_i,d}^*$  then ▷ by Condition (5) in Theorem 3
9:         ▷ These values can also be recorded to save computation
10:        Divide  $r_i$  into  $r_{i1} = [a_{i1}, b_{i1}]$  and  $r_{i2} = [a_{i2}, b_{i1}]$  along the max gap (Figure 1).
11:         $\Pr(r_{i1}) = \Pr(r_i) \frac{|P(r_{i1})|}{n_i}$ ,  $\Pr(r_{i2}) = \Pr(r_i) - \Pr(r_{i1})$ 
12:         $\mathcal{B}' = \mathcal{B}' \cup \{r_{i1}, r_{i2}\}$ 
13:      else  $\mathcal{B}' = \mathcal{B}' \cup \{r_i\}$ 
14:    if  $\mathcal{B}' \neq \mathcal{B}$  then  $\mathcal{B} = \mathcal{B}'$ 
15:    else return  $\mathcal{B}, \Pr(\cdot)$ 

```

Remark 2.1. Zero probability is not desirable in some applications; it can be avoided by adding pseudo count (Laplace smoother) α in line 11, i.e., $\Pr(r_{i1}) = \Pr(r_i) \frac{|P(r_{i1})| + \alpha}{n_i + 2\alpha}$. Density $d(r_i)$ is recovered by $\Pr(r_i) / \prod_{j=1}^d (b_{ij} - a_{ij})$.

Remark 2.2. The binary tree shown in Figure 1 can be constructed as a byproduct and the user can specify the deepest level to terminate the algorithm.

and apply Theorem 1 to $\tilde{f}(\tilde{x})$, we have

Proof. The proof is in Proposition 10 of Section 8 in [17]. \square

$$\left| \int_{[0,1]^d} \tilde{f}(\tilde{x}) d\tilde{x} - \frac{1}{n} \sum_{i=1}^n \tilde{f}(\tilde{x}_i) \right| \leq D_{n_i}^*(\tilde{P}) V_{HK}^{[0,1]^d}(\tilde{f})$$

Lemma 3.3. Let

$$D_{n,d}^* = \inf_{x_1, \dots, x_n \in [0,1]^d} D_n^*(x_1, \dots, x_n)$$

we have

$$D_{n,d}^* \leq cd^{1/2} n^{-1/2}$$

for all $n, d = 1, 2, \dots$, with a multiplicative constant c .

Remark 3.1. It is also shown that $c \leq 10$ in [2]. The asymptotic behavior of the star discrepancy on n is much better (e.g., Halton sequence [16] has $D_n^* = O((\log n)^d/n)$); but it does not necessarily mean that the uniform bound which is valid for all d and n cannot be of order $n^{-1/2}$.

Proof. The proof is quite technical and presented in Theorem 3 of [10]. \square

Theorem 2. f is defined on d -dimensional hyper-rectangle $[a, b]$ and $P = \{x_1, \dots, x_n \in [a, b]\}$. Then we have

$$\begin{aligned} \left| \int_{[a,b]} f(x) dx - \frac{\prod_{i=1}^d (b_i - a_i)}{n} \sum_{i=1}^n f(x_i) \right| \\ \leq \prod_{i=1}^d (b_i - a_i) D_n^*(\tilde{P}) V_{HK}^{[a,b]}(f) \end{aligned} \quad (4)$$

where $\tilde{P} = \{\tilde{x}_i = (\frac{x_{i1}-a_1}{b_1}, \dots, \frac{x_{id}-a_d}{b_d})\}_{i=1}^n$

Proof. Define $\tilde{f}(\tilde{x}) = f(x)$, where $\tilde{x} = (\frac{x_1-a_1}{b_1}, \dots, \frac{x_d-a_d}{b_d})$

From Lemma 3.2, $V_{HK}^{[0,1]^d}(\tilde{f}) = V_{HK}^{[a,b]}(f)$; $\int_{[0,1]^d} \tilde{f}(\tilde{x}) d\tilde{x} = (\prod_{i=1}^d (b_i - a_i))^{-1} \int_{[a,b]} f(x) dx$ by change of variables and $\tilde{f}(\tilde{x}_i) = f(x_i)$ by definition. Hence, (4) follows immediately. \square

We are ready to state our main theorem,

Theorem 3. f is defined on hyper-rectangle $[0, 1]^d$ with $V_{HK}^{[0,1]^d}(f) < \infty$ and the sub-rectangles $\{[a_i, b_i]\}_{i=1}^l$ are a split of $[0, 1]^d$. Let $x_1, \dots, x_N \in [0, 1]^d$ be an iid sample set drawn from distribution $p(x)$ defined on $[0, 1]^d$ and $P_i = \{x_{i1}, \dots, x_{in_i}, n_i \in \mathbb{N}^+\}$ are points in each sub-region. Consider a piecewise constant density estimator

$$\hat{p}(x) = \sum_{i=1}^l d_i \mathbf{1}\{x \in [a_i, b_i]\}$$

where $d_i = (\prod_{j=1}^d (b_{ij} - a_{ij}))^{-1} n_i / N$, i.e., the empirical probability. In each sub-region $[a_i, b_i]$, P_i satisfies

$$D_{n_i}^*(\tilde{P}_i) \leq \alpha_i D_{n_i,d}^* \quad (5)$$

where $\alpha_i = \sqrt{\frac{N}{n_i d}} \frac{\theta}{c}$ and θ is a positive constant; \tilde{P}_i is defined as $\{\tilde{x}_j = (\frac{x_{j1}-a_{i1}}{b_{i1}}, \dots, \frac{x_{jd}-a_{id}}{b_{id}})\}_{j=1}^{n_i}$ then

$$\left| \int_{[0,1]^d} f(x) \hat{p}(x) dx - \frac{1}{N} \sum_{i=1}^N f(x_i) \right| \leq \frac{\theta}{\sqrt{N}} V_{HK}^{[0,1]^d}(f) \quad (6)$$

Remark 3.2. α_i controls the “relative” uniformity of the points and is adapted for P_i , i.e., it imposes more restricted constraint on the region containing large proportion of the sample (n_i/N).

Remark 3.3. In Monte Carlo methods, the convergence rate of $\frac{1}{N} \sum_{i=1}^N f(x_i)$ to $\int_{[0,1]^d} f(x)p(x)dx$ is of order $O(1/\sqrt{N})$ as long as variance of $f(x)$ under $p(x)$ is bounded; our density estimate is optimal in the sense that it achieves the same rate of convergence.

Remark 3.4. There are many other $\hat{p}(x)$ satisfying (6) such as the empirical distribution in the extreme or kernel density estimation with sufficiently small bandwidth. Our density estimation is attractive in the sense that it provides a very sparse summary of the data whereas captures the landscape of the underlying distribution; moreover, the piecewise constant function does not suffer from many local bumps as kernel density estimation does.

Proof. Apply Theorem 2 to each $[a_i, b_i], i = 1, \dots, l$, we have

$$\begin{aligned} & \left| \int_{[a_i, b_i]} f(x)dx - \frac{\prod_{j=1}^d (b_{ij} - a_{ij})}{n_i} \sum_{j=1}^{n_i} f(x_{ij}) \right| \\ & \leq \prod_{j=1}^d (b_{ij} - a_{ij}) D_{n_i}^*(\tilde{P}_i) V_{HK}^{[a_i, b_i]}(f) \end{aligned} \quad (7)$$

and by triangular inequality, we have

$$\begin{aligned} & \left| \int_{[0,1]^d} f(x)\hat{p}(x)dx - \frac{1}{N} \sum_{i=1}^N f(x_i) \right| \\ & \leq \sum_{i=1}^l d_i \left| \int_{[a_i, b_i]} f(x)dx - \frac{1}{d_i N} \sum_{j=1}^{n_i} f(x_{ij}) \right| \end{aligned} \quad (8)$$

By the definition of d_i , $d_i N = (\prod_{j=1}^d (b_{ij} - a_{ij}))^{-1} n_i$; combine with Theorem 2, (5), (7) and Lemma 3.3, we have

$$\begin{aligned} & \sum_{i=1}^l d_i \left| \int_{[a_i, b_i]} f(x)dx - \frac{1}{d_i N} \sum_{j=1}^{n_i} f(x_{ij}) \right| \\ & \leq \sum_{i=1}^l d_i \prod_{j=1}^d (b_{ij} - a_{ij}) D_{n_i}^*(\tilde{P}_i) V_{HK}^{[a_i, b_i]}(f) \\ & \leq \sum_{i=1}^l \frac{n_i}{N} \sqrt{\frac{N}{n_i d}} \frac{\theta}{c} D_{n_i, d}^* V_{HK}^{[a_i, b_i]}(f) \\ & \leq \sum_{i=1}^l \frac{n_i}{N} \sqrt{\frac{N}{n_i d}} \frac{\theta}{c} c d^{1/2} n_i^{-1/2} V_{HK}^{[a_i, b_i]}(f) \\ & \text{Apply Lemma 3.1} \\ & = \frac{\theta}{\sqrt{N}} \sum_{i=1}^l V_{HK}^{[a_i, b_i]}(f) = \frac{\theta}{\sqrt{N}} V_{HK}^{[0,1]^d}(f) \end{aligned}$$

Corollary 3.1. For any hyper-rectangle $A = [a, b] \subset (0, 1)^d$. Let $\hat{P}(A) = \int_A \hat{p}(x)dx$ and $P(A) = \int_A p(x)dx$, then $|\hat{P}(A) - P(A)|$ converges to 0 at order $O(1/\sqrt{N})$ uniformly.

Remark 3.5. The total variation distance between probability measures \hat{P} and P is defined via $\delta(\hat{P}, P) = \sup_{A \in \mathcal{B}} |\hat{P}(A) - P(A)|$, where \mathcal{B} is the Borel σ algebra of $[0, 1]^d$; in contrast, Corollary 3.1 restricts A to be rectangles.

Proof. In Monte Carlo methods, the convergence rate of $\frac{1}{N} \sum_{i=1}^N f(x_i)$ is of order $O(\frac{\text{std}(f)}{\sqrt{N}})$. Let $f(x) = \mathbf{I}\{x \in [a, b]\} = \mathbf{I}_{[a, b]}$ be defined on $[0, 1]^d$, we have $\text{var}(f) = P(A)(1 - P(A)) \leq 1/4$; thus, this error is bounded uniformly.

If another indicator function \tilde{f} is defined on $[\tilde{a}, \tilde{b}] \subset (0, 1)^d$, then let

$$\begin{aligned} \phi_j(\tilde{x}_j) = & \frac{a_j}{\tilde{a}_j} \tilde{x}_j \mathbf{I}_{[0, \tilde{a}_j]} + (a_j + \frac{b_j - a_j}{\tilde{b}_j - \tilde{a}_j} (\tilde{x}_j - \tilde{a}_j)) \mathbf{I}_{[\tilde{a}_j, \tilde{b}_j]} \\ & + (b_j + \frac{1 - b_j}{1 - \tilde{b}_j} (\tilde{x}_j - \tilde{b}_j)) \mathbf{I}_{[\tilde{b}_j, 1]} \end{aligned}$$

and $\phi(\tilde{x}) = \prod_{j=1}^d \phi_j(\tilde{x}_j)$ and apply Lemma 3.2, we have $V_{HK}^{[0,1]^d}(\tilde{f}) = V_{HK}^{[0,1]^d}(f)$; thus, the left term of (6) is bounded uniformly.

Combining the two parts, the theorem follows by triangular inequality. \square

4 COMPUTATIONAL ASPECTS

There are no explicit formulas for calculating $D_n^*(x_1, \dots, x_n)$ and $D_{n, d}^*$ except for low dimensions. If we replace α_i in (5) and apply Lemma 3.3, we actually intend to control $D_n^*(\tilde{P}_i)$ by $\theta\sqrt{N}/n_i$. There are several ways to approximate $D_n^*(x_1, \dots, x_n)$: 1) E. Thiérmard presents an algorithm to compute the star discrepancy within a user specified error by partitioning the unit cube into subintervals [20], [21], [6]; 2) A genetic algorithm to calculate the lower bounds is proposed in [19] and is used in our experiments; 3) A new randomized algorithm based on threshold accepting is developed in [8]. Comprehensive numerical tests indicate that it improves on other algorithms, especially in higher dimension $20 \leq d \leq 50$. The interested readers are referred to the original articles for implementation details.

In dealing with large data, several simple observations can be exploited to save computation: 1) it is trivial that $\max_{j=1, \dots, d} D_n^*({x_{ij}}_{i=1}^n) \leq D_n^*({x_i}_{i=1}^n)$. Let $x_{(i)j}$ be the i th smallest element in $\{x_{ij}\}_{i=1}^n$, then $D_n^*({x_{ij}}_{i=1}^n) = \frac{1}{2n} + \max_{i=1}^n |x_{(i)j} - \frac{2i-1}{2n}|$ [5], which has complexity $O(n \log n)$. Hence $\max_{j=1, \dots, d} D_n^*({x_{ij}}_{i=1}^n)$ can be used to compare against $\theta\sqrt{N}/n$ first before calculating $D_n^*({x_i}_{i=1}^n)$; 2) $\theta\sqrt{N}/n$ is large when n is small, but $D_n^*({x_i}_{i=1}^n)$ is bounded above by 1; 3) $\theta\sqrt{N}/n$ is tiny when n is large and $D_n^*({x_i}_{i=1}^n)$ is bounded below by $c_d \log^{(d-1)/2} n^{-1}$ with some constant c_d depending on d [7]; thus we can keep splitting without checking (5) when

\square

$\theta\sqrt{N}/n \leq \epsilon$, where ϵ is a small positive constant (say 0.001) specified by the user. This strategy may introduce few more sub-rectangles, but the running time gain is considerable.

Another approximation works well in practice is by replacing star discrepancy with computationally attractive \mathcal{L}_2 star discrepancy, i.e., $D_n^{(2)}(x_1, \dots, x_n) = (\int_{[0,1]^d} |\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{x_i \in [0,a]} - \prod_{i=1}^d a_i|^2 da)^{1/2}$; in fact, several statistics to test uniformity hypothesis based on $D_n^{(2)}$ are proposed in [13]; however, the theoretical guarantee in Theorem 3 is no longer valid. By Warnock's formula [5],

$$[D_n^{(2)}(x_1, \dots, x_n)]^2 = \frac{1}{3^d} - \frac{2^{1-d}}{n} \sum_{i=1}^n \prod_{k=1}^d (1 - x_{ik}^2) + \frac{1}{n^2} \sum_{i,j=1}^n \prod_{k=1}^d \min\{1 - x_{ik}, 1 - x_{jk}\}$$

$D_n^{(2)}$ can be computed in $O(n \log^{d-1} n)$ by K. Frank and S. Heinrich's algorithm [5]. At each scan of \mathcal{B} in Algorithm 1, the total complexity is at most $\sum_{i=1}^l O(n_i \log^{d-1} n_i) \leq \sum_{i=1}^l O(n_i \log^{d-1} n) \leq O(n \log^{d-1} n)$.

5 EXPERIMENTAL RESULTS

5.1 Simulation

1) To demonstrate the methods and visualize the results, we simulate our algorithm through 3 2-dimensional data sets generated from 3 distributions respectively, i.e.,

$$x \sim \mathcal{N}(\mu, \Sigma) \mathbf{1}\{x \in [0, 1]^2\}$$

with

$$\mu = \begin{pmatrix} 0.50 \\ 0.50 \end{pmatrix}, \Sigma = \begin{pmatrix} 0.08, 0.02 \\ 0.02, 0.02 \end{pmatrix}$$

and

$$x \sim \frac{1}{2} \left(\mathcal{N}(\mu_1, \Sigma_1) + \mathcal{N}(\mu_2, \Sigma_2) \right) \mathbf{1}\{x \in [0, 1]^2\}$$

with

$$\mu_1 = \begin{pmatrix} 0.50 \\ 0.25 \end{pmatrix}, \Sigma_1 = \begin{pmatrix} 0.04, 0.01 \\ 0.01, 0.01 \end{pmatrix}$$

$$\mu_2 = \begin{pmatrix} 0.50 \\ 0.75 \end{pmatrix}, \Sigma_2 = \begin{pmatrix} 0.04, 0.01 \\ 0.01, 0.01 \end{pmatrix}$$

and

$$x \sim \frac{1}{3} \left(\beta_{2,5} \beta_{5,2} + \beta_{4,2} \beta_{2,4} + \beta_{1,3} \beta_{3,1} \right)$$

where \mathcal{N} is the Gaussian distribution and β is the beta distribution. The size of each data set is 10,000. As shown in the first row of Figure 2, we draw the partitions on 2D and render the estimated densities with a color map; the corresponding contours of true densities are embedded for comparison purpose.

2) To evaluate the theoretical bound (6), we choose 3 simple reference functions with dimension $d = 2, 5$ and 10 respectively, i.e., $f_1(x) = \sum_{i=1}^n \sum_{j=1}^d x_{ij}^{1/2}$, $f_2(x) = \sum_{i=1}^n \sum_{j=1}^d x_{ij}$, $f_3(x) = (\sum_{i=1}^n \sum_{j=1}^d x_{ij}^{1/2})^2$ and samples

are generated from

$$p(x) = \frac{1}{2} \left(\prod_{i=1}^d \beta_{15,5}(x_i) + \prod_{i=1}^d \beta_{5,15}(x_i) \right)$$

The error $|\int_{[0,1]^d} f_i(x)p(x)dx - \int_{[0,1]^d} f_i(x)\hat{p}(x)dx|$ is bounded by

$$\left| \int_{[0,1]^d} f_i(x)p(x)dx - \frac{1}{n} \sum_{j=1}^n f_i(x_j) \right| + \left| \int_{[0,1]^d} f_i(x)\hat{p}(x)dx - \frac{1}{n} \sum_{j=1}^n f_i(x_j) \right|$$

where $\hat{p}(x)$ is the estimated density; the first term is controlled by $O(n^{-1/2})$ which is well known in Monte Carlo methods and the second term is controlled by (6), thus the error is of order $O(n^{-1/2})$. By varying the data size, the relative error v.s. sample size is plotted on log-log scale for each dimension in the second row of Figure 2, their standard errors are obtained through generating 10 replicas under same distributions. Interestingly, the linear pattern shows up as expected.

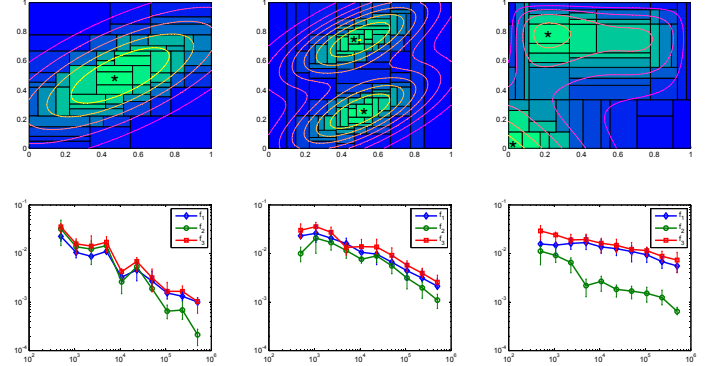


Fig. 2. First row: simulation on 2D data with 3 different densities; the modes are marked by stars. Second row: simulation on 2, 5 and 10 dimensional data (from left to right) with sample functions f_1, f_2, f_3 .

5.2 Mode Detection

A direct application of the piecewise constant density is to detect modes [4], i.e., the dense areas or local maxima on the domain. The modes of our density estimator is defined as

Definition 5.1. A mode of the piecewise constant density is a sub-rectangle in the partition that its density is largest among all its neighbors as indicated by the stars in Figure 2.

In order to compare our method with OPT or BSP¹ in terms of running times and performance

1. The source codes are obtained from the authors. Their implementation language is C++; in contrast, our method is implemented in Matlab. For small data, the latency of Matlab dominates the computing time as shown in the first block of Table 2.

in mode detection, we simulate samples from $x \sim (\sum_{i=1}^4 \pi_i \mathcal{N}_i(\mu_i, \Sigma)) \mathbf{1}\{x \in [0, 1]^d\}$ with $d = \{2, 3, 4, 5, 6\}$ and $n = \{10^3, 10^4, 10^5\}$ respectively, where

$$\begin{pmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \\ \mu_4 \end{pmatrix} = \begin{pmatrix} 1/4 & 1/4 & 1/2 & \cdots & 1/2 \\ 1/4 & 3/4 & 1/2 & \cdots & 1/2 \\ 3/4 & 1/4 & 1/2 & \cdots & 1/2 \\ 3/4 & 3/4 & 1/2 & \cdots & 1/2 \end{pmatrix}_{4 \times d}$$

and $\Sigma = 0.01\mathbf{I}$, i.e., the identity matrix, $\pi = (1/4, 1/4, 1/4, 1/4)$. The system where the comparison is performed is Ubuntu 13.04, amd64 8-core Opteron 2384 (SB X6240, 0916FM400J); 31.42GB RAM, 10GB swap.

The results of mode detection are summarized in Table 1 and running time are in Table 2, the standard error is obtained by generating 20 replicas for each (d, n) pair. It is shown that density estimates by OPT and BSP have more modes generally. One possible explanation is that MAPs of OPT and BSP try to find the global optimizer of all possible binary partitions, they tend to overfit the data and result in a partition with many noisy sub-regions; in contrast, DED makes myopic decisions and the possible choices for splitting are limited, but one can still bound the overall integration error by controlling the discrepancy adaptively as (5).

5.2.1 Flow Cytometry

Flow cytometry allows to measure simultaneously multiple characteristics of a large number of cells and is a ubiquitous and indispensable technology in medical settings. One effort of current research is to identify homogeneous sub-populations of cells automatically instead of manual gating, which is criticized for its subjectivity and non-scalability. There are a large amount of recent literatures concerning on auto gating and clustering, see [1] and many references therein.

In order to apply our method, we regard each cell as one observation in the sample space, i.e., if there are n markers attached to a single cell, then the whole data set is generated from a hypothetical n dimensional distribution. Mature cell populations concentrate in some high density areas, which can be easily identified in the binary partitioned space by Definition 5.1.

One practical issue needs to be addressed for most of the Cytometry analysis techniques: there is asymmetry in sub-populations; by optimizing a predefined loss function, it is possible that some sparse yet crucial populations are overlooked if the algorithms take most of the efforts to control the loss in denser areas. A remedy for this issue is to perform a down-sampling [1], [18] step to roughly equalize the densities among populations then up-sampling after populations are identified; however, this step is dangerous that it may fails to sample enough cells in sparse populations, as a result, these populations are lost in the down-sampled data. In contrast, our approach does not require down-sampling step, and the asymmetry among populations are captured by their densities.

For the mouse bone marrow data studied in [18], we choose the 8 markers (SSA-C, CD11b, B220, TCR- β , CD4, CD8, c-kit, Sca-1) that are relevant to the cell types of interests; the number of cells is $\sim 380,000$ after removing mutli-cell aggregates and co-incident events. 13 sub-populations are identified by our algorithm ([18] and its supplementary materials), the results are summarized in Figure 3.

5.2.2 Image Segmentation

Following [14] in which a new density estimation via histogram transforms is proposed, we conduct a similar experiment dealing with color image segmentation. The author in [14] reports that the result of his new algorithm is “barely the same” as that of others, thus we use mean shift with Gaussian kernel density estimator as the benchmark, which is publicly available in the GUI version of Edge Detection and Image SegmentatiON (EDISON) system [4]. For each pixel, we concatenate its LUV feature space representation with its coordinates to form a 5-dim *joint domain* [4] representation. Our method are used to learn a 5-dim piecewise constant density. After identifying the modes according to Definition 5.1, we use k -means to group the pixels with the metric

$$d(x_1, x_2) = (\|x_1^r - x_2^r\|_2^2 + \lambda \|x_1^s - x_2^s\|_2^2)^{\frac{1}{2}}$$

we write $x = (x^r, x^s)$ corresponding to the *range* (color) domain and *spatial* domain; λ controls the relative importance of spatial difference, for example, a large λ tends to connect adjacent pixels even if their colors are very different. Each cluster obtained from k -means corresponds to several patches in the original image; and each pixel is replaced by the average color in the patch it belongs to.

Once each pixel is process as above, some region connecting or pruning algorithms are employed to eliminate spurious patches. For easy of comparison, we employ the APIs in EDISON system to merge patches with its default parameters. The images are chosen from USC-SIPI Image Database [22] and are rescaled to 256×256 pixels by bicubic interpolation. The results are summarized in Figure 4.

5.3 Some Other Applications

Density Topology Exploration and Visualization. The connectivity graph (DG) or level set tree [24] is widely used to represent energy landscapes of systems; it summarizes the hierarchy among various local maxima and minima in the configuration space; its topology is a tree and each inner node on the tree is a changing point that merges two or more independent regions in the domain. With the density estimation at hand, one may construct DG for samples instead of a given energy or density function. Unlike KDE that suffers from many local bumps and results in an overly complicated DG, (3) is well suited for this purpose, partially because it smoothes out the minor fluctuations and takes only

TABLE 1
The average number of modes detected by OPT, BSP and DED for each pair (d, n) respectively.

	#modes($n = 10^3$)			#modes($n = 10^4$)			#modes($n = 10^5$)		
	OPT	BSP	DED	OPT	BSP	DED	OPT	BSP	DED
d									
2	5.1(1.1)	4.4(1.2)	3.8(0.4)	5.9(1.3)	5.7(0.9)	4.0(0.7)	8.1(3.1)	7.1(2.3)	4.8(0.8)
3	3.2(0.5)	3.7(0.8)	2.4(0.5)	4.7(1.2)	6.2(1.7)	3.5(0.4)	7.7(2.9)	6.9(1.3)	4.4(0.5)
4	4.1(0.8)	4.6(1.0)	2.7(0.4)	6.1(2.1)	5.7(1.8)	3.0(0.9)	6.4(2.0)	7.2(3.3)	4.2(0.4)
5	3.3(0.6)	4.1(1.5)	2.1(0.6)	6.6(1.7)	7.8(2.2)	3.7(0.8)	8.7(2.0)	8.1(3.2)	4.2(1.1)
6	4.7(1.2)	4.3(1.4)	3.1(0.5)	5.9(1.9)	7.5(2.9)	4.2(1.0)	9.1(1.7)	8.2(4.4)	5.1(1.3)

TABLE 2
The average running time of OPT, BSP and DED for each pair (d, n) respectively. The stars indicates that the running time exceeds 3600s (in order to save space).

	#modes($n = 10^3$)			#modes($n = 10^4$)			#modes($n = 10^5$)		
	OPT	BSP	DED	OPT	BSP	DED	OPT	BSP	DED
d									
2	0.4(0.0)	1.2(0.1)	1.7(0.1)	2.8(0.1)	23.2(6.4)	11.2(0.9)	42.9(0.3)	263.1(44.9)	95.8(3.6)
3	0.8(0.0)	1.6(0.3)	2.2(0.4)	13.3(0.1)	27.7(8.4)	17.1(1.9)	252(2.8)	422.8(91.7)	143.7(2.0)
4	1.7(0.1)	3.5(0.2)	3.3(0.8)	137.7(10.2)	42.3(5.3)	22.6(1.8)	*	684.3(80.0)	192.4(5.1)
5	75.6(3.3)	4.9(0.3)	3.2(0.7)	1731.7(17.7)	138.2(9.7)	21.3(2.2)	*	1547.9(155.6)	231.6(6.8)
6	251.3(7.9)	5.1(0.4)	3.8(0.7)	*	179.1(13.4)	30.0(2.1)	*	*	285.4(10.2)

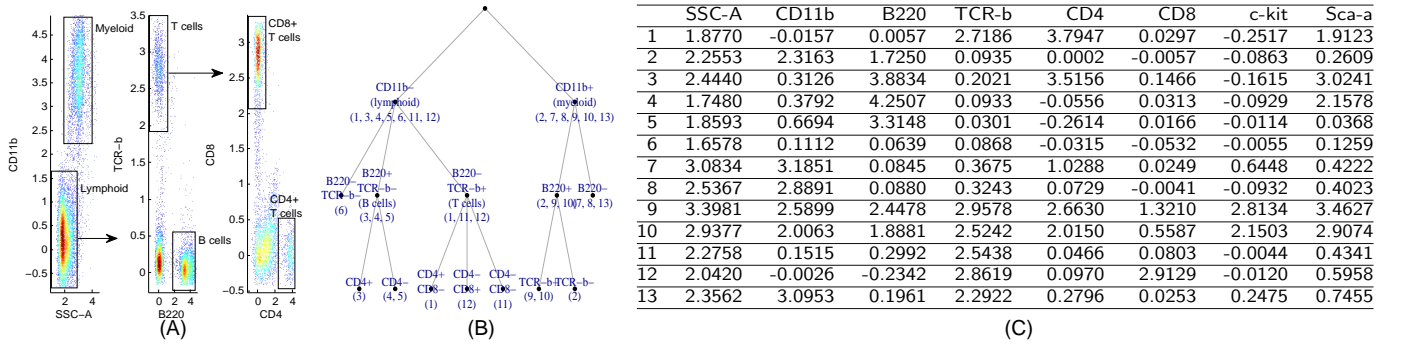


Fig. 3. (A): an illustrative gating sequence, the cell type in each gate is attached; (B) there are 13 modes detected by our algorithm, and we arrange these modes into a hierarchical dendrogram: at first level, they are grouped by expression levels of CD11b; subsequently, the CD11b- modes are grouped according to B220 and TCR-b then further splitted according to CD4 and CD8 on the next level; the CD11b+ modes are grouped by B220 then by TCR-b; (C) the details of the expression levels of each mode.

limited number of values; moreover, the simple structure of (3) makes the construction of such graph easy (i.e., one can just scan through each r_i in decreasing order of $d(r_i)$). The DG of (3) not only reveals the modes of the density on its leaves, it also provides a tool to visualize high dimensional data hierarchically; for example, in fiber tractography [11], DG is used to visualize and analyze topography in fiber streamlines interactively.

We demonstrate that how our piecewise density function can be used to construct level set trees in Figure 5. The basic pipeline is to scan sub-rectangles sequentially

according to the decreasing order of their densities and agglomerate the sub-rectangles according to their adjacency.

Multi-level Feature Extraction. The density of each observation is available after learning the density function (3) and each sub-rectangle groups the observations with similar densities. These densities contains important non-linearity within the data which is hard to capture by standard transformations. We can augment the feature space of the sample by appending their corresponding densities. Through varying the deepest levels

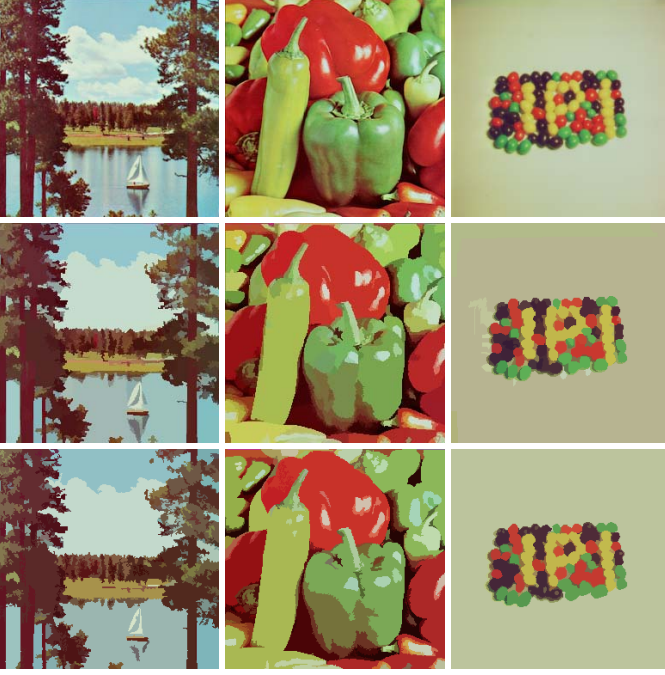


Fig. 4. a) 1st row: original images; 2nd row: segmentation by mean shift with Gaussian kernel with default parameters; 3rd row: segmentation by DED. b) 1st column: lake; 2nd column: pepper; 3rd column: beans.

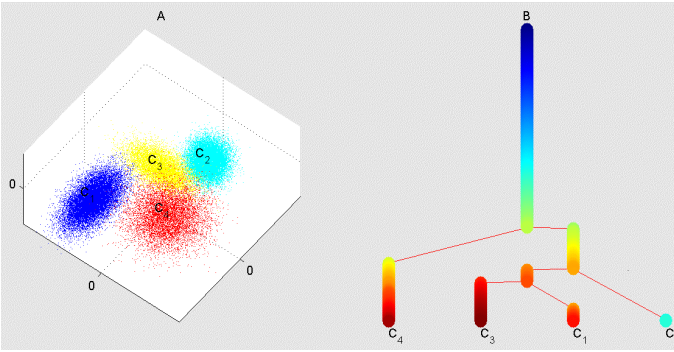


Fig. 5. Left (A): the samples are generated from a Gaussian Mixture with 4 modes. Right (B): the level set tree. The clusters are annotated by C_1, C_2, C_3, C_4 .

(Remark 2.2), the densities learned from different levels are included in the features; specifically, let $\hat{p}_{l_1}, \dots, \hat{p}_{l_k}$ are learned densities of sample $\{x_i\}_{i=1}^n$ by controlling the deepest levels to be l_1, \dots, l_k respectively, then the learned features are

$$\{(\hat{p}_{l_1}(x_i), \dots, \hat{p}_{l_k}(x_i))\}_{i=1}^n$$

This multi-level feature extraction technique has potential applications in representation learning.

6 CONCLUSION AND FUTURE WORK

We have developed a nonparametric density estimation framework based on discrepancy criteria, proved its theoretical properties and shown that it is applicable to different types of problems. We point out several future research directions of interest: 1) current approach deals with continuous features, but how to extend our theories and algorithm to handle (unordered) categorical data? 2) coordinate wise partition limits the approximation capability, recent progress [3] in Quasi Monte Carlo on simplex provides us a possible alternative to use more flexible partition schemes. 3) theoretically, how to sharpen Corollary 3.1 in order to enlarge the class of Borel sets rather than rectangles or more aggressively, to prove the consistency? 4) a throughout comparison is necessary to understand the empirical differences between our method and OPT or BSP.

ACKNOWLEDGMENTS

Kun Yang is supported by General Wang Yaowu Stanford Graduate Fellowship and The Simons Math+X fellowship; Wing Hung Wong is supported by NSF grants DMS 0906044 and 1330132.

REFERENCES

- [1] N. Aghaeepour, G. Finak, H. Hoos, T. R. Mosmann, R. Brinkman, R. Gottardo, R. H. Scheuermann, F. Consortium, and D. Consortium. Critical assessment of automated flow cytometry data analysis techniques. *Nature methods*, 2013.
- [2] C. Aistleitner. Covering numbers, dyadic chaining and discrepancy. *Journal of Complexity*, 27(6):531–540, 2011.
- [3] K. Basu and A. B. Owen. Low discrepancy constructions in the triangle. *arXiv preprint arXiv:1403.2649*, 2014.
- [4] D. Comaniciu and P. Meer. Mean shift: A robust approach toward feature space analysis. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 24(5):603–619, 2002.
- [5] C. Doerr, M. Gnewuch, and M. Wahlström. Calculation of discrepancy measures and applications. *Preprint*, 2013.
- [6] M. Gnewuch. Bracketing numbers for axis-parallel boxes and applications to geometric discrepancy. *Journal of Complexity*, 24(2):154–172, 2008.
- [7] M. Gnewuch. Entropy, randomization, derandomization, and discrepancy. In *Monte Carlo and quasi-Monte Carlo methods 2010*, pages 43–78. Springer, 2012.
- [8] M. Gnewuch, M. Wahlström, and C. Winzen. A new randomized algorithm to approximate the star discrepancy based on threshold accepting. *SIAM Journal on Numerical Analysis*, 50(2):781–807, 2012.
- [9] L. Györfi. *Principles of nonparametric learning*, volume 434. Springer, 2002.
- [10] S. Heinrich, E. Novak, G. W. Wasilkowski, and H. Wozniakowski. The inverse of the star-discrepancy depends linearly on the dimension. *ACTA ARITHMETICA-WARSZAWA-*, 96(3):279–302, 2000.
- [11] B. P. Kent, A. Rinaldo, F.-C. Yeh, and T. Verstynen. Mapping topographic structure in white matter pathways with level set trees. *arXiv preprint arXiv:1311.5312*, 2013.
- [12] L. Kuipers and H. Niederreiter. *Uniform distribution of sequences*. Courier Dover Publications, 2012.
- [13] J.-J. Liang, K.-T. Fang, F. Hickernell, and R. Li. Testing multivariate uniformity and its applications. *Mathematics of Computation*, 70(233):337–355, 2001.
- [14] E. López-Rubio. A histogram transform for probability density function estimation. *IEEE transactions on pattern analysis and machine intelligence*, 2013.
- [15] L. Lu, H. Jiang, and W. H. Wong. Multivariate density estimation by bayesian sequential partitioning. *Journal of the American Statistical Association*, 108(504):1402–1410, 2013.

- [16] A. B. Owen. Quasi-monte carlo sampling. *Monte Carlo Ray Tracing: Siggraph*, pages 69–88, 2003.
- [17] A. B. Owen. Multidimensional variation for quasi-monte carlo. In *International Conference on Statistics in honour of Professor Kai-Tai Fang's 65th birthday*, pages 49–74, 2005.
- [18] P. Qiu, E. F. Simonds, S. C. Bendall, K. D. Gibbs Jr, R. V. Bruggner, M. D. Linderman, K. Sachs, G. P. Nolan, and S. K. Plevritis. Extracting a cellular hierarchy from high-dimensional cytometry data with spade. *Nature biotechnology*, 29(10):886–891, 2011.
- [19] M. Shah. A genetic algorithm approach to estimate lower bounds of the star discrepancy. *Monte Carlo Methods and Applications*, 16(3-4):379–398, 2010.
- [20] E. Thiérmard. Computing bounds for the star discrepancy. *Computing*, 65(2):169–186, 2000.
- [21] E. Thiérmard. An algorithm to compute bounds for the star discrepancy. *journal of complexity*, 17(4):850–880, 2001.
- [22] A. Weber. The usc-sipi image database. *Signal and Image Processing Institute of the University of Southern California*. URL: <http://sipi.usc.edu/services/database>, 1997.
- [23] W. H. Wong and L. Ma. Optional pólya tree and bayesian inference. *The Annals of Statistics*, 38(3):1433–1459, 2010.
- [24] Q. Zhou and W. H. Wong. Energy landscape of a spin-glass model: Exploration and characterization. *Physical Review E*, 79(5):051117, 2009.